# Performance-Related Activities at NERSC/CRD

**Charlene Yang**

**Application Performance Specialist**
**NERSC, LBNL**

# Activities at NERSC/CRD

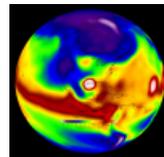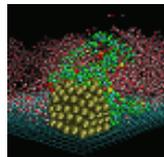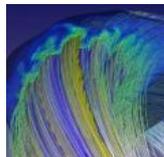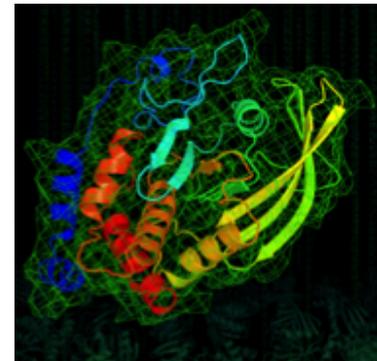- **Roofline performance model**
  - NESAP, vendor integration
  - Performance portability, scaling trajectories
  - Instruction Roofline, integer Roofline, mixed Precision
  - ERT, energy Roofline, Roofline for FPGA

- **LDMS for mass performance data collection**
  - *#SBATCH --profile=<tool>:<group>*
  - <tool> = vtune, likwid, ldms; <group> = flops, mem, bandwidth, …

- **PAPI for Roofline**
  - *#SBATCH --profile=timemory:roofline*

# Roofline Performance Model

# Roofline Performance Model

- **Sustainable performance is bound by**

$$\text{GFLOP/s} = \min \begin{cases} \textbf{Peak GFLOP/s} \\ \textbf{AI * Peak GB/s} \end{cases}$$

- **Arithmetic Intensity (AI) =**

  **FLOPs / Bytes**

- **How did this come about?**
  - **→ A CPU DRAM example**
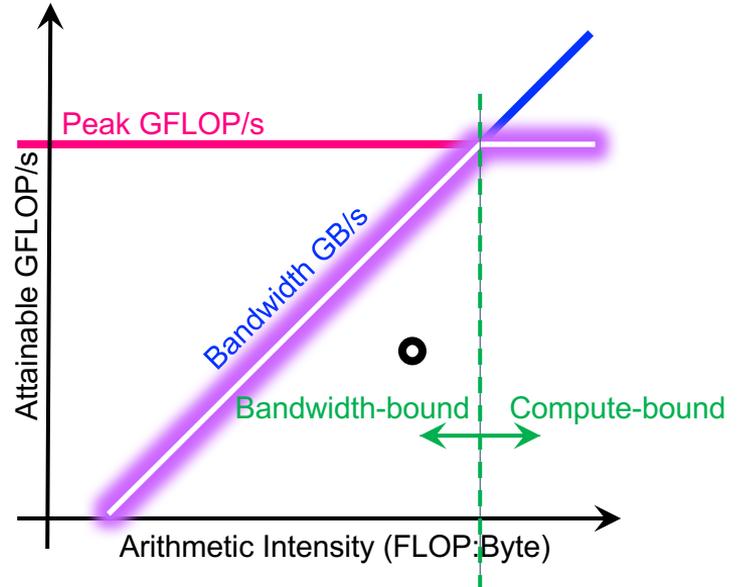


Peak GFLOP/s

Bandwidth GB/s

Attainable GFLOP/s

Arithmetic Intensity (FLOP:Byte)

Bandwidth-bound    Compute-bound
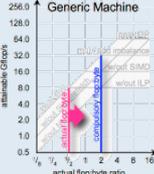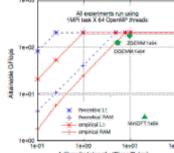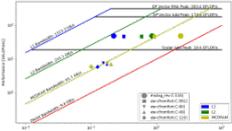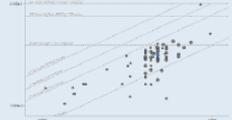
Transition @ AI ==
Peak GFLOP/s / Peak GB/s ==
'Machine Balance'

# The Roofline Chronical

|  | 2005 - 2011 | 2013 - 2016 | 2017 - 2019 | Future |
|---|---|---|---|---|
| **Research** | ▪ Developed **foundations** for the Roofline Model<br>▪ Applied to kernels using canonical flops and bytes  |  | ▪ Developed performance counter Rooflines for **CPUs and GPUs**<br>▪ Roofline for **Simulations** and **Machine Learning**<br>▪ Incorporated VPU%, divides, integer operations | ▪ FPGAs, CGRAs, AI processors, …<br>▪ Asymmetric memory hierarchies<br>▪ Horizontal data movement<br>▪ Effects of **extreme heterogeneity** |
| **Prototype** |  | ▪ Created the **ERT** prototype for CPUs and GPUs<br>▪ Quantified CUDA UVM effects  | ▪ Collaboration with CRD, Intel and NVIDIA on **hierarchical Roofline**  | ▪ **Integer/instruction/non-FP** Rooflines<br>▪ Rooflines that serialize data transfers (vs. assume overlap)<br>▪ Integration with compilers/runtimes |
| **Production** |  |  | ▪ Roofline model incorporated into **Intel Advisor**<br>▪ Installed at NERSC, LANL, *etc*  | ▪ **Roofline for GPUs (multiple vendors)**<br>▪ Roofline for FPGAs/CGRAs<br>▪ Integer/instruction/non-FP Rooflines<br>▪ CISC/DL instructions |

# The Roofline People

## Researchers…

- Sam Williams (Roofline Lead, LBL/CRD)
- Doug Doefler (LBL/NERSC)
- Khaled Ibrahim (LBL/CRD)
- Nan Ding (LBL/CRD)
- Yunsong Wang (LBL/NERSC)
- Jack Deslippe (LBL/NERSC)
- Lenny Oliker (RAPIDS deputy, LBL/CRD)
- Terry Ligocki (LBL/CRD)
- Brian Van Straalen (LBL/CRD)
- Aleksandar Ilic (INESC, Portugal)
- Diogo Marques (INESC, Portugal)

## Vendors/Industry…

- Zakhar Matveev (Intel)
- Max Katz, Magnus Strengert (NVIDIA)
- Constantios Evangelinos (IBM)
- Protonu Basu (Facebook; formerly LBL/CRD)
- Linda Lo (Facebook; formerly U. Utah)
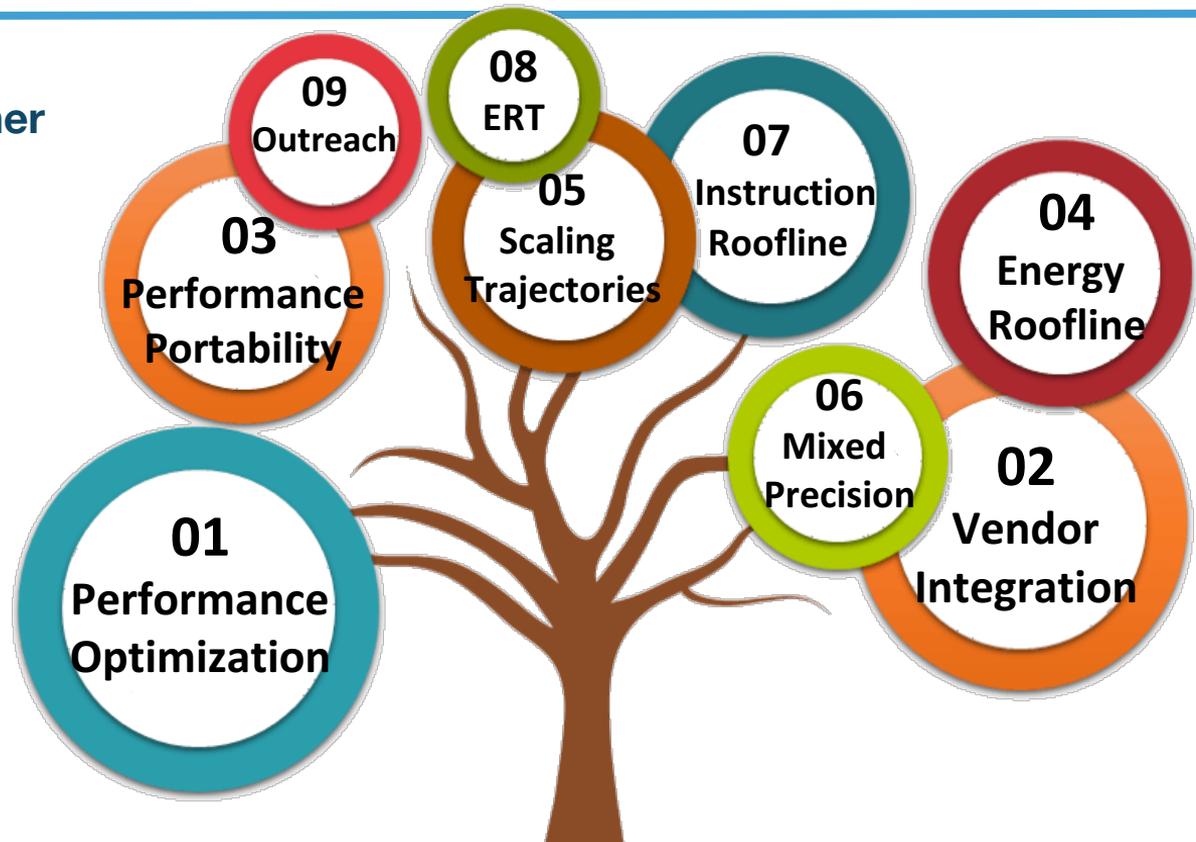- David Patterson (Google, formerly UC Berkeley)
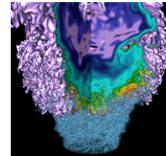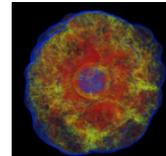
*Thank You!*

# The Roofline Tree

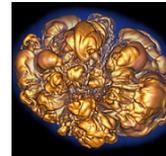**Brings People Together**
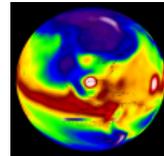
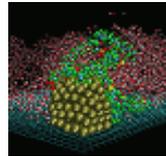- **NESAP**
- **CRD**
- **Intel**
- **NVIDIA**
- **all HPCers**



**Roofline Performance Model**

# 1. Roofline drives optimization

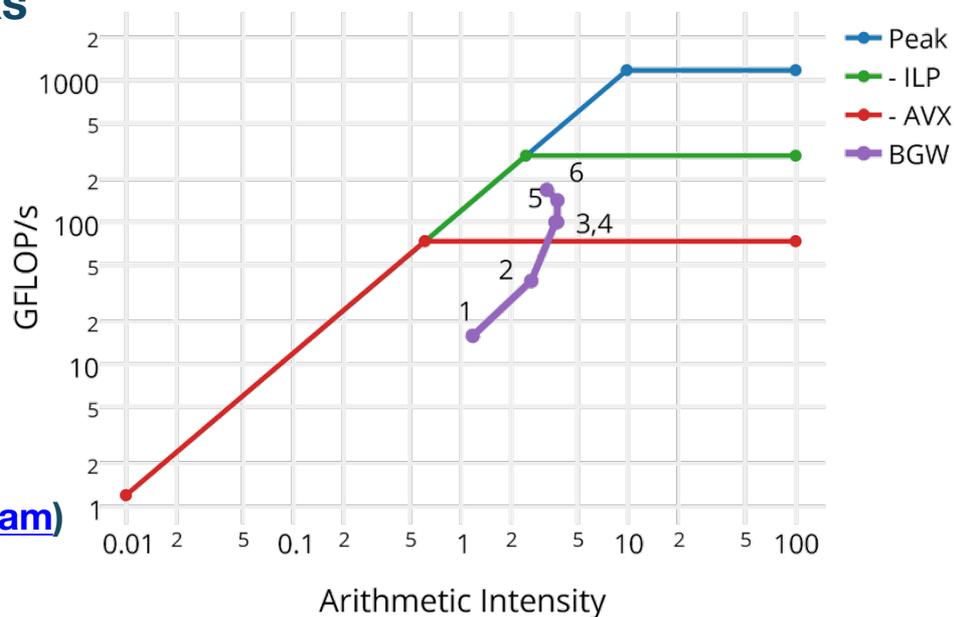## NESAP

# Roofline Drives Optimization

**The Roofline Model**

- helps you identify the bottlenecks
- guides you through optimization
- tells you when to stop

**An example:**

- **NESAP for Cori - BerkeleyGW**

(**NERSC Exascale Scientific Application Program**)
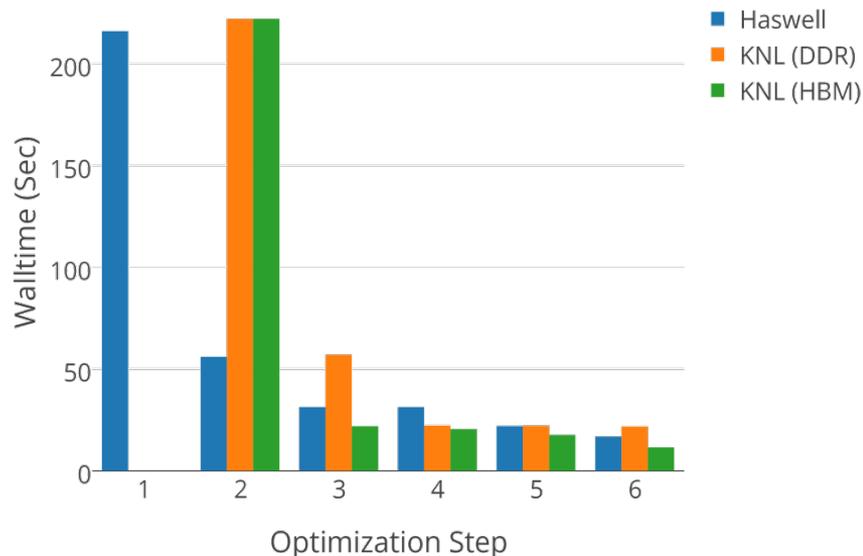


Haswell Roofline Optimization Path

# Roofline Drives Optimization

**Optimization Path for Kernel-C (Sigma):**

1. **Add OpenMP**
2. **Initial Vectorization**
   - **loop reordering**
   - **conditional removal**
3. **Cache-Blocking**
4. **Improved Vectorization**
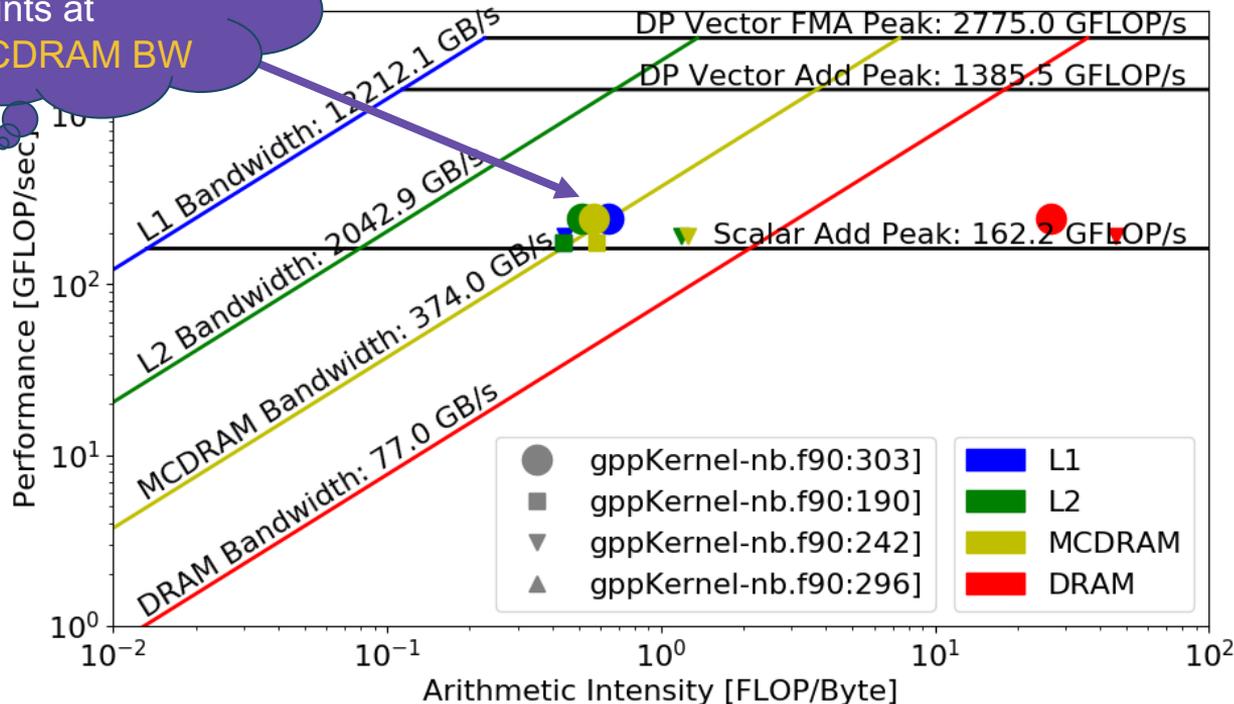   - **divides**
5. **Hyper-threading**



Sigma Optimization Process

# Example 1: GPP, KNL, Cache Blocking
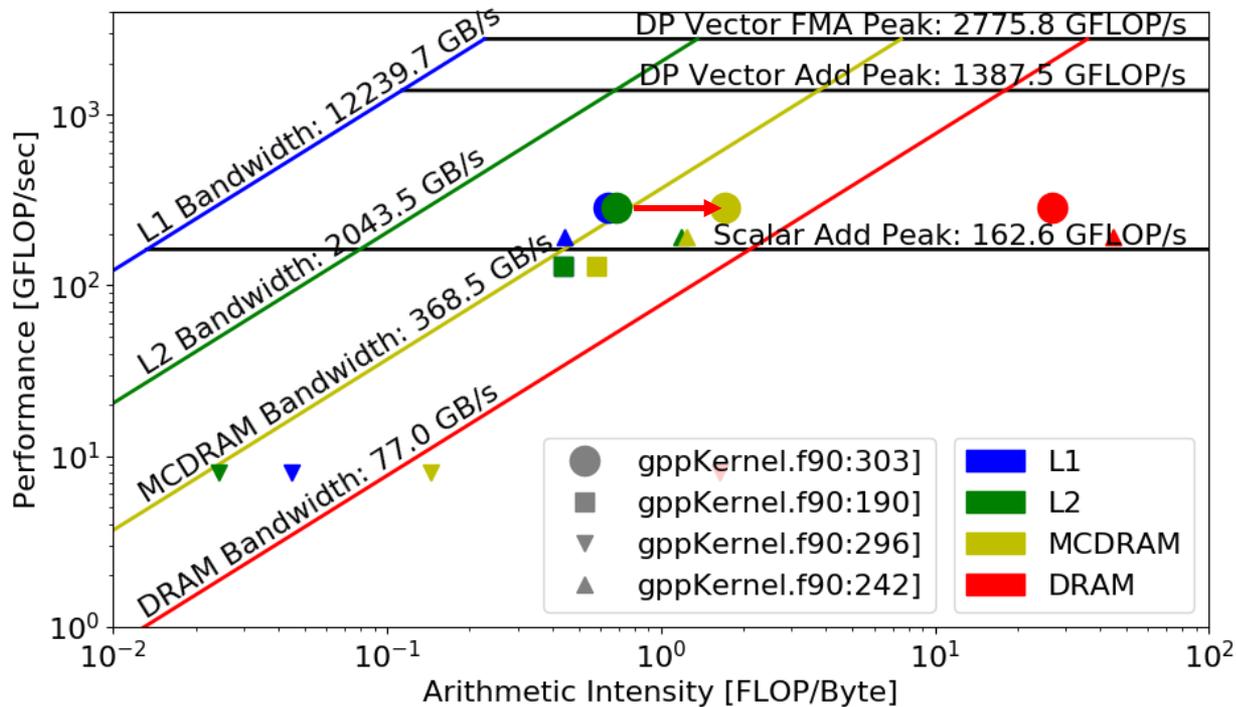


Overlapping points at MCDRAM BW

242 GFflop/s, **Bound by MCDRAM Bandwidth**

Most Flops in the main loop (O)

Read/Write 2MB of data per inner loop iteration ➤ No reuse of data in L1/L2, shown by overlapping points at MCDRAM bandwidth

BW Bound ➤ Increase MCDRAM AI by adding cache locality

# Example 1: GPP, KNL, Cache Blocking



Cache blocking implemented to achieve L2 data reuse

**3x Increase in MCDRAM AI**

Performance increased from 242 to 287 GFlop/s (+18%)

Why not 3x Flops increase?
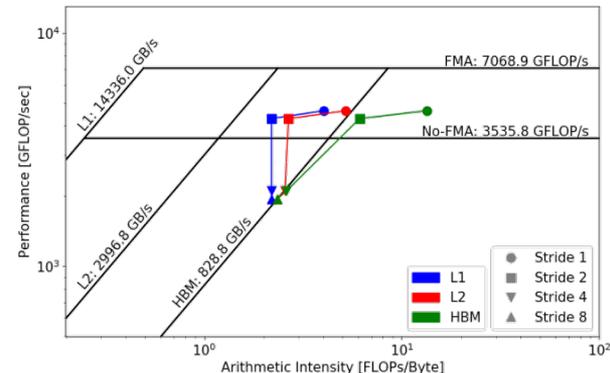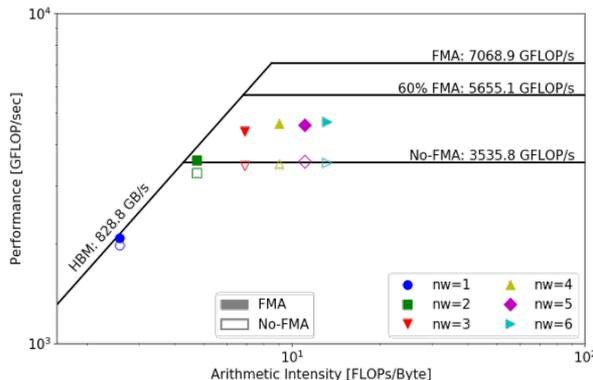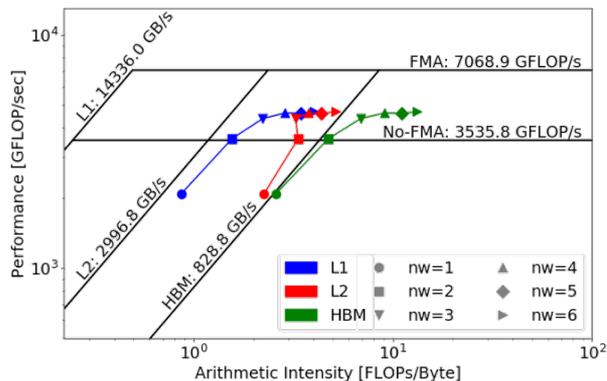➤ Not BW bound any more, divide, shuffle and unpack instructions involved

• T. Koskela, Z. Matveev, C. Yang, A. Adetokunbo, R. Belenov, P. Thierry, Z. Zhao, R. Gayatri, H. Shan, L. Oliker, J. Deslippe, R. Green, and S. Williams, A Novel Multi-Level Integrated Roofline Model Approach for Performance Characterization, ISC'2018 Research Paper, Jun 24-28 2018, Frankfurt
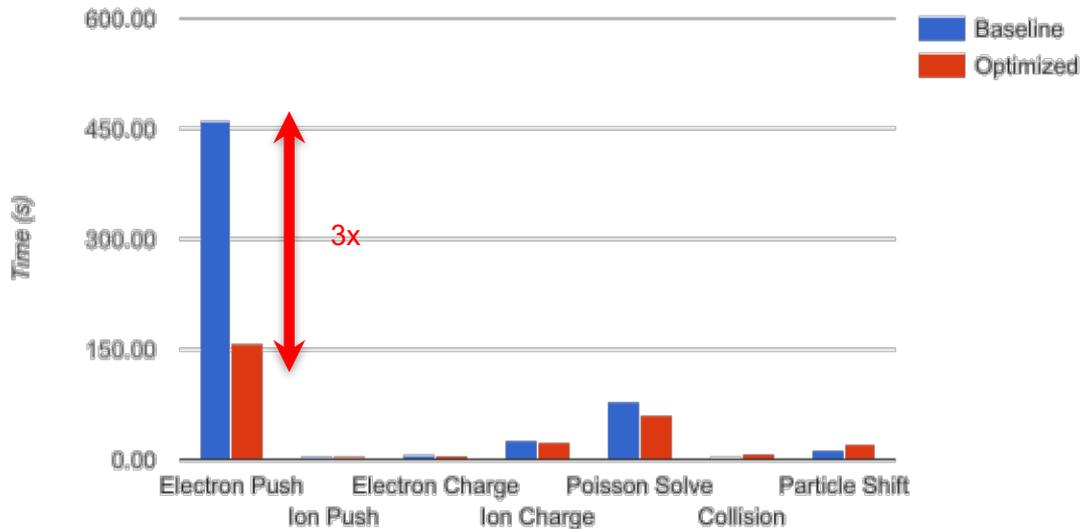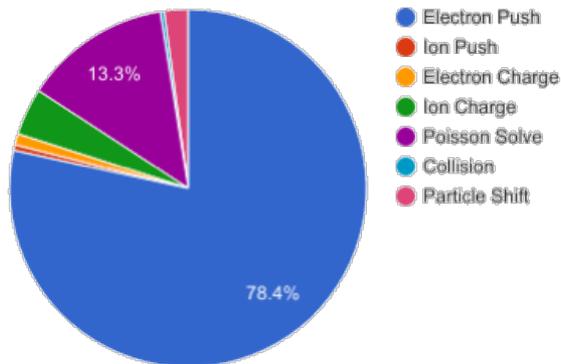
**Three experiments to study the effects of**

- **cache reuse (varying `nw` from 1 to 6)**

- **instruction mix (FMA vs. Mul/Add)**

- **memory coalescing**

```
do band = 1, nbands       #blockIdx.x
   do igp = 1, ngpown      #blockIdx.y
      do ig = 1, ncouls   #threadIdx.x
         do iw = 1, nw     #unrolled
            compute; reductions
```
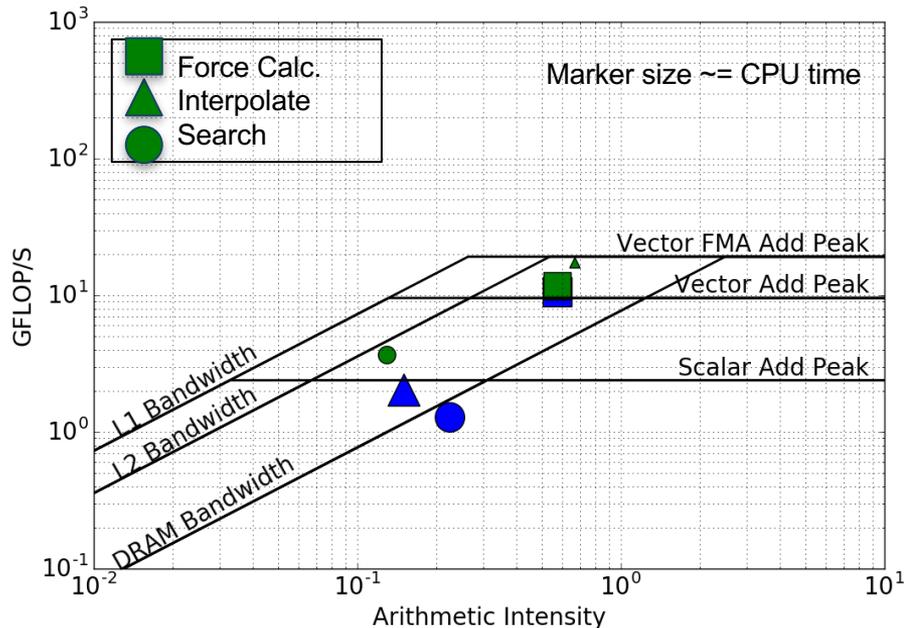
# Example 2: XGC1, KNL



(Left) Hotspots for unoptimized XGC1 on 1024 Cori KNL nodes in Quad-Flat mode;
(Right) Speedup in XGC1 Electron Push routine after back porting the optimizations made in ToyPush kernel

# Example 2: ToyPush from XGC1



- **Force Kernel:**
- close to vector add peak
- not much optimization done

- **Interpolate Kernel:**
- L1 blocking, indirect memory access
- memory alignment, more efficient vectorization
- **10x speedup**, closer to vector FMA peak

- **Search Kernel:**
- multiple exits, simd private, enable vectorization
- **3x speedup**, closer to L2 bandwidth roof
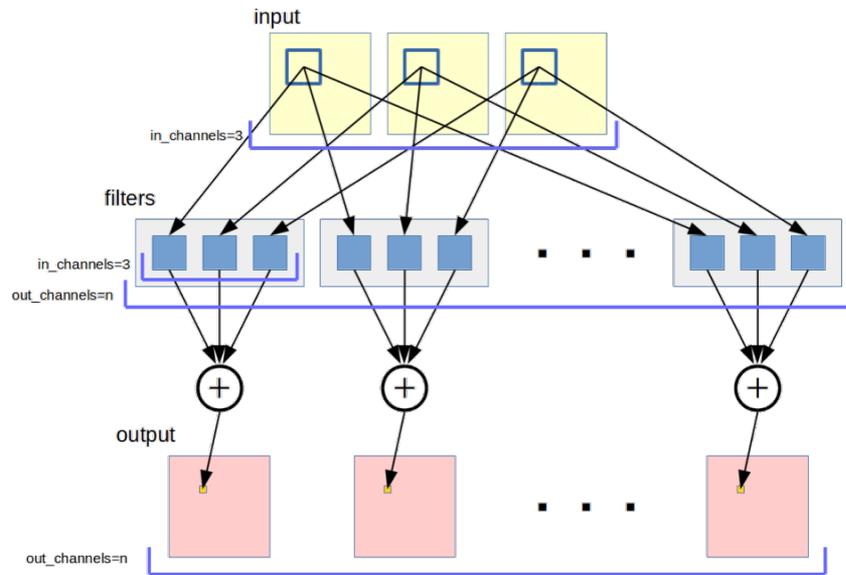
- Code is available at
- https://github.com/tkoskela/toypush

# Example 3: conv2d from TensorFlow



- **Kernel tf.nn.conv2d**

https://www.tensorflow.org

$$B_{nhwc} = \sum_{m=0}^{C-1} \sum_{k_h=0}^{K_H-1} \sum_{k_w=0}^{K_W-1} A_{n\ h+k_h\ w+w_h\ m}\ K_{k_h\ k_w\ m\ c}$$
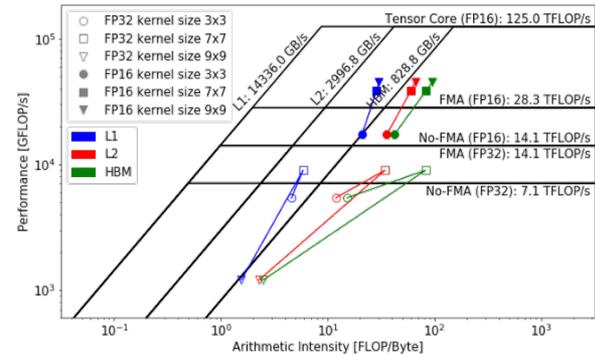
# Example 3: TF / Forward Pass



## #Batch Size

- Constant performance(**no!**)
- **FP16 performance anti-correlated with batch size**
- Performance << TC peak
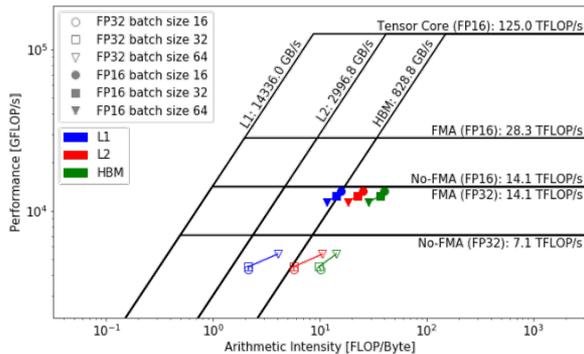- Transformation kernels
- Low L2 locality

## #Filters

- Intensity ∝ #Filters
- Low L2 data locality
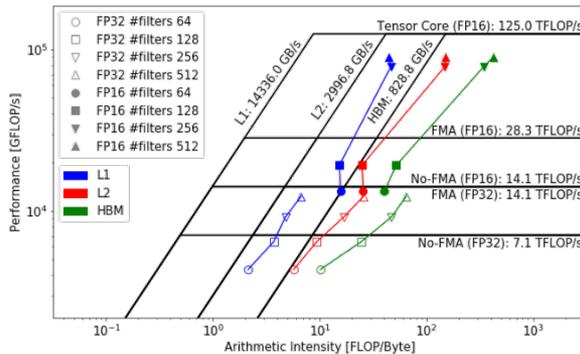- Some use of TC's (>FP16 FMA)… *partial TC ceiling*

## #Kernel Size

- Intensity ∝ kernel size
- Low L2 data locality
- **Autotuner switched FP32 algorithm to FFT at 9x9**
- Some use of TC's (>FP16 FMA)… *partial TC ceiling*

# Example 3: TF / Backward Pass

## #Batch Size

o Autotuner chose different (**better**) algorithm for FP32 with batch size = 64 (boost)

## #Filters

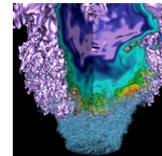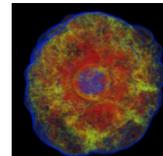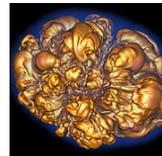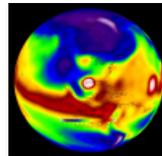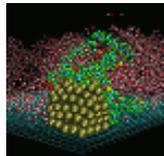o Close to FP16 TC peak
o Close to FP32 FMA peak

## #Kernel Size

o Good FP32 performance trend (almost peak)

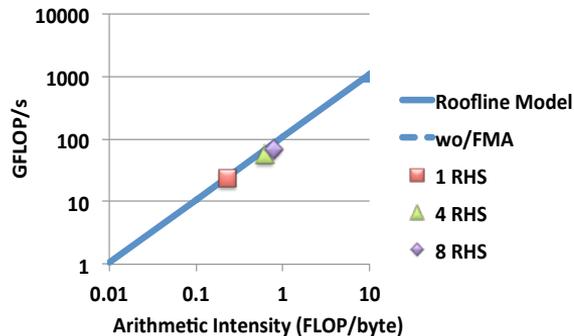o **Autotuner chose to run 9x9 FP16 in FP32 !!**
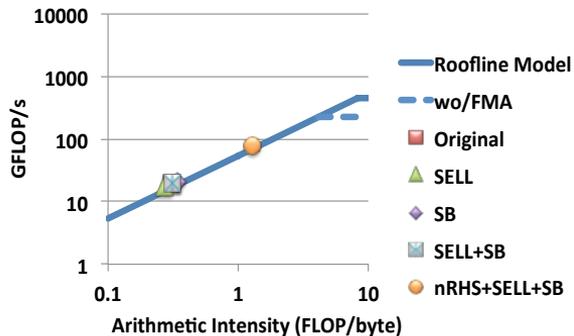
# 2. Vendor Integration

**Intel VTune, LIKWID, Intel Advisor, NVIDIA nvprof**
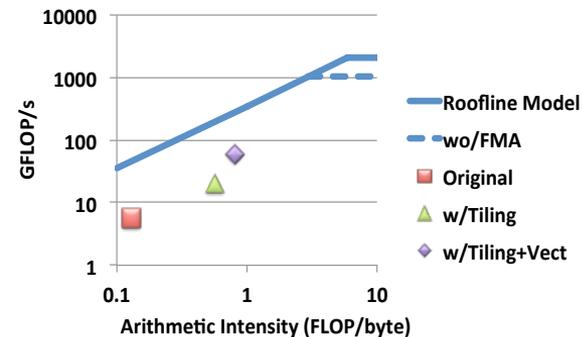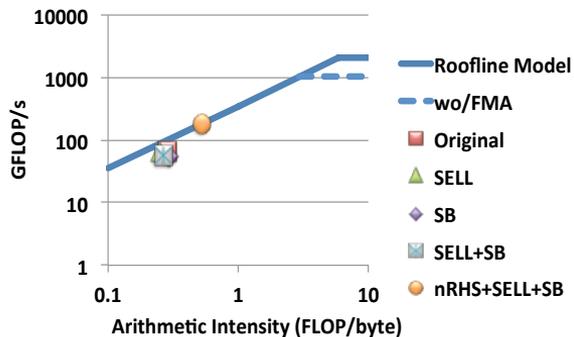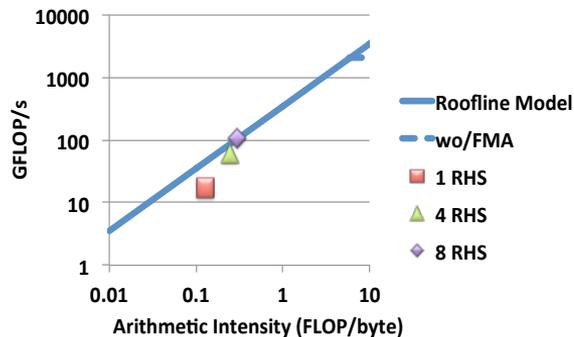
# Data Collection on Intel CPUs

**Way 1:**

- **Intel SDE** for FLOPs  (emulation)

- **Intel VTune** for DRAM bytes (HW counters)

- **Runtime**


- **DRAM Roofline** only


- **Used by NESAP for Cori**

  - **NERSC Exascale Science Application Program**

  - http://www.nersc.gov/users/application-performance/measuring-arithmetic-intensity/

# Data Collection on Intel CPUs



**DRAM Rooflines of NESAP Codes**
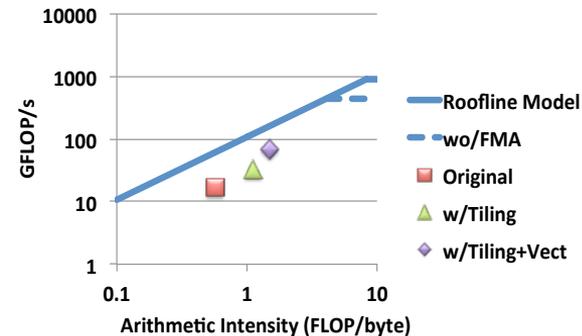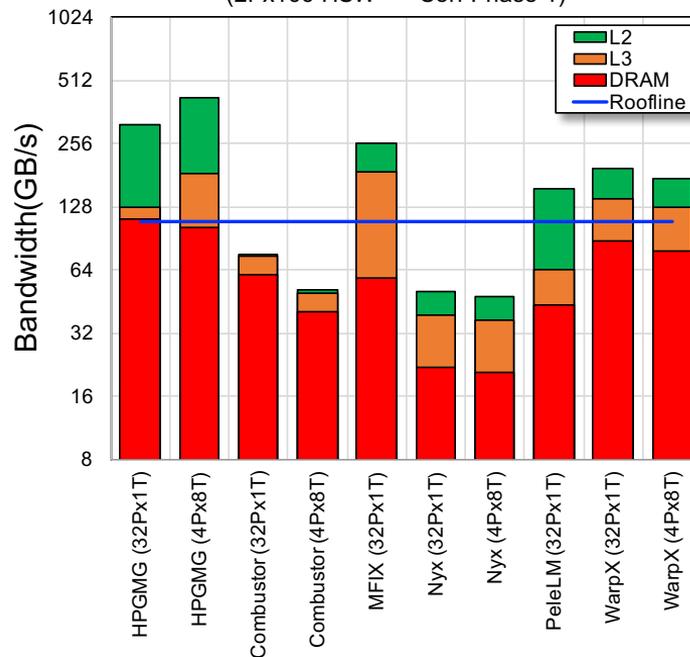
# Data Collection on Intel CPUs

**Way 2:**

- **LIKWID** for FLOPs and bytes
  - Both are based on HW counters
- **Runtime**

- **Hierarchical** Roofline

- **Limited by quality of HW counters**
- **High-level characterization, no callstack (need instrumentation)**



AMReX Application Characterization
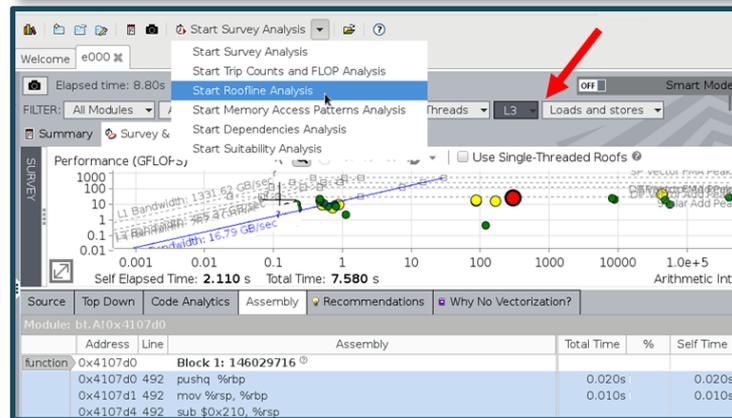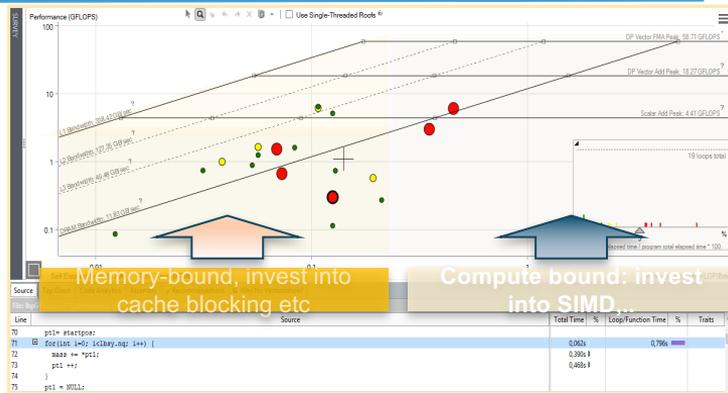(2Px16c HSW == Cori Phase 1)

https://github.com/RRZE-HPC/likwid

# Data Collection on Intel CPUs

**Way 3:**

- **Intel Advisor, Roofline feature**
- **Instrument applications automatically**
  - **one dot per loop nest/function**
- **FLOPs, bytes and runtime**


- **Hierarchical Roofline**


- **Integrates with other Advisor capabilities**
- **Benchmarks target system**

# Data Collection on Intel CPUs

## New features in Intel Advisor 2019

**(picture courtesy of Z. Matveev)**

https://software.intel.com/en-us/intel-advisor-2019-release-notes
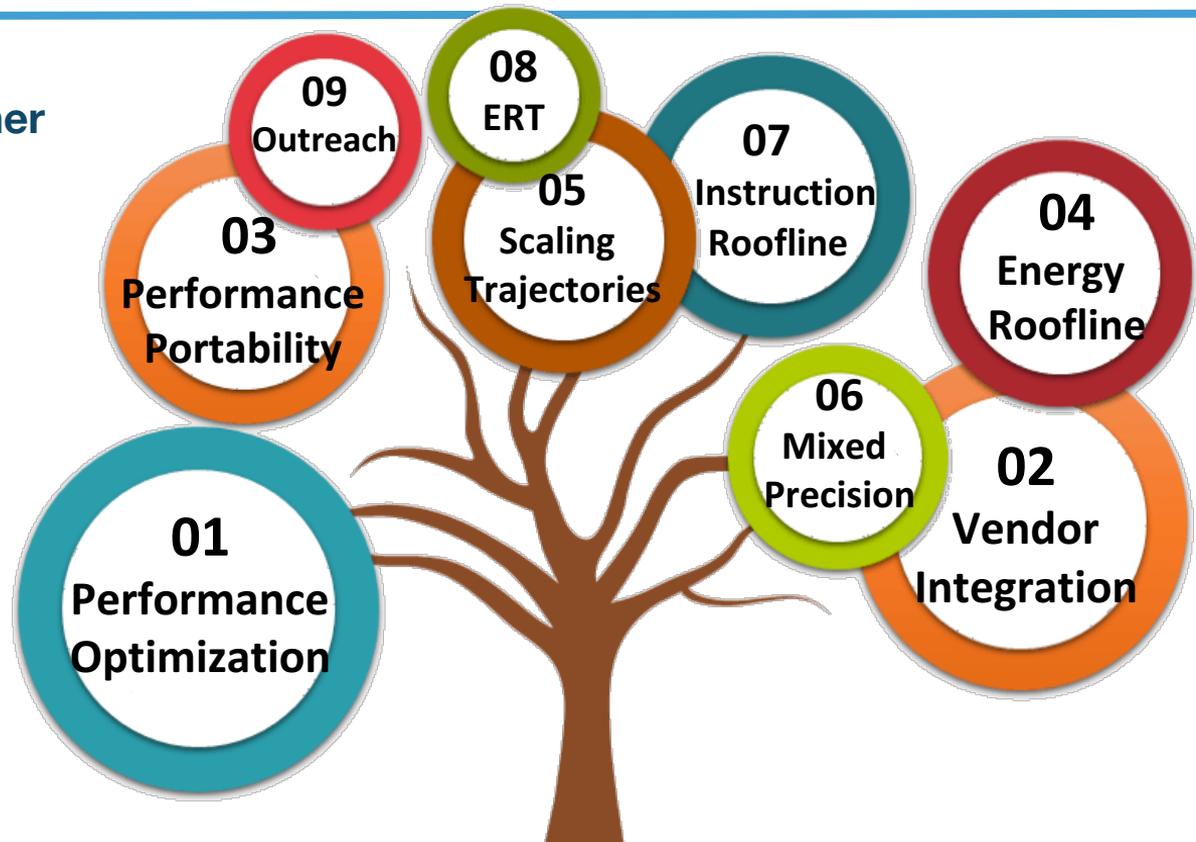
# Data Collection on NVIDIA GPUs

- **Still manual at this stage, but we have a recipe using nvprof.**

- **Runtime:**
  - **Internal timers or `nvprof --print-gpu-trace`**
- **FLOPs:**
  - **DP/SP/HP counters and metrics, `nvprof --metrics` `'flop_count_dp/sp/hp'` or `` `tensor_precision_fu_utilization' ``**
- **Bytes for different cache levels:**
  - **Bytes = (read transactions + write transactions) x transaction size**
  - **`nvprof --metrics` `'metric_name'` e.g. `gld/gst_transactions`**

- **Hierarchical Roofline**

- C. Yang, S. Williams, Hierarchical Roofline Analysis for GPUs: Accelerating Performance Optimization for the NERSC-9 Perlmutter System, CUG'2019, May 5-9 2019, Montreal, Canada

# The Rest of the Tree

**Brings People Together**

- **NESAP**
- **CRD**
- **Intel**
- **NVIDIA**
- **all HPCers**



**Roofline Performance Model**

# 3. Performance Portability

**Definition, Metric, Roofline, KNL, V100**

# Introduction

- No consensus on the definition or metric for performance portability

- But Pennycook *et al*…

$$\Phi(a, p, \boldsymbol{H}) = \begin{cases} \dfrac{|\boldsymbol{H}|}{\sum_{i \in \boldsymbol{H}} \dfrac{1}{e_i(a,p)}} & \text{if } i \text{ is supported, } \forall i \in \boldsymbol{H} \\ 0 & \text{otherwise} \end{cases}$$

- Architectural Efficiency [Williams *et al*]

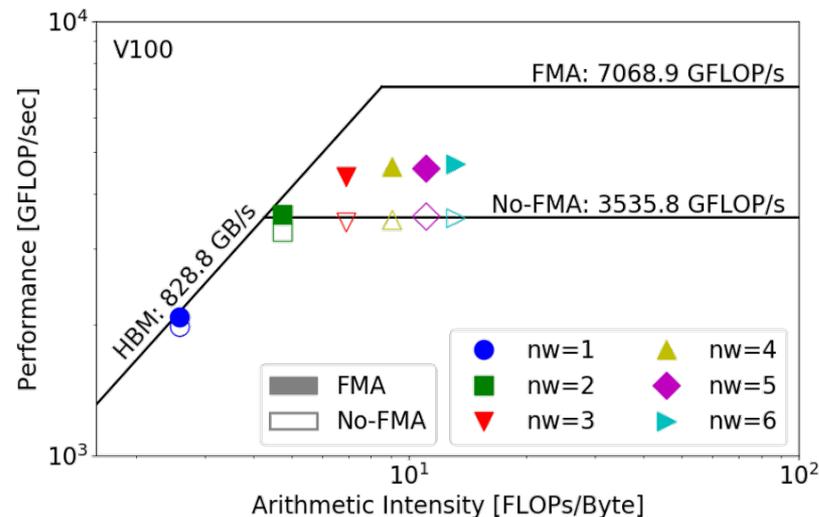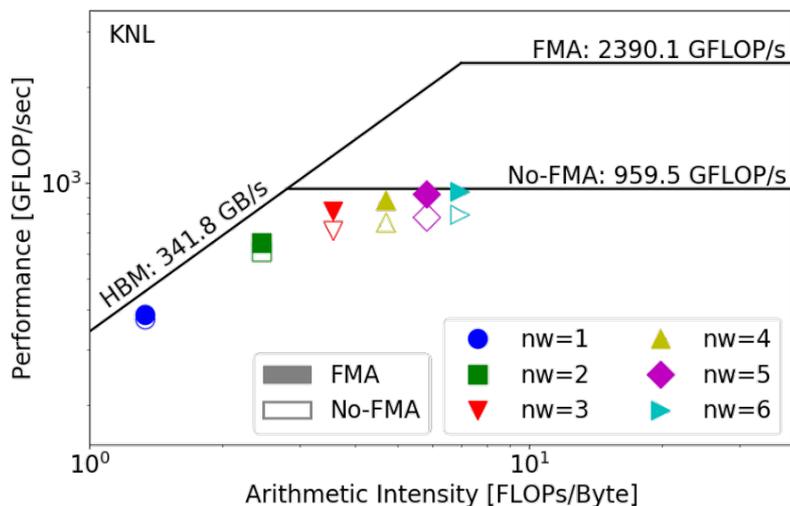$$e_i(a, p) = \frac{P_i(a, p)}{\min(F_i, \ B_i \times I_i(a, p))}$$

**Peak FLOP/s**

**Arithmetic Intensity**

**Peak Bandwidth**

# Bottleneck Changes

- **Bottleneck shifts at $nw = 2$ on KNL vs. V100 (no-FMA performance)**
- **Easier to achieve no-FMA ceiling on V100 than KNL, due to higher ratio of instruction issue bandwidth vs. instruction execution bandwidth**

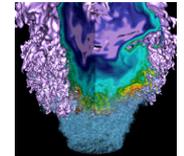- C. Yang, R. Gayatri, T. Kurth, P. Basu, Z. Ronaghi, A. Adetokunbo, B. Friesen, B. Cook, D. Doerfler, L. Oliker, J. Deslippe, S. Williams, An Empirical Roofline Methodology for Quantitatively Assessing Performance Portability, SC'2018 P3HPC Workshop, Nov 11-16 2018, Dallas
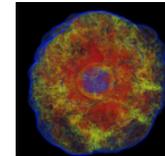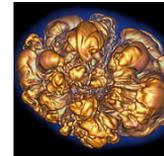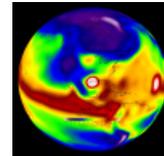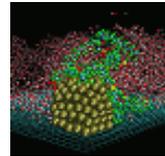
# Bottleneck Changes

- **No FMA: performance portability consistently > 80%**
- **FMA: benefit is far less than 2x at high $nw$; architectural efficiency suffers (so does performance portability)**
- **Could regain some architectural efficiency if non-floating-point vector operations were considered**

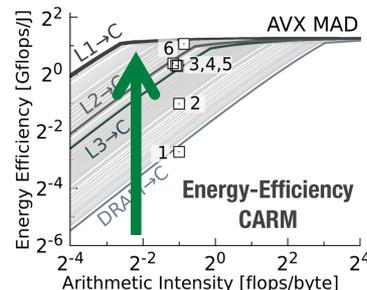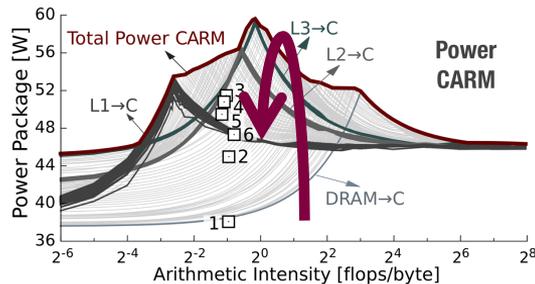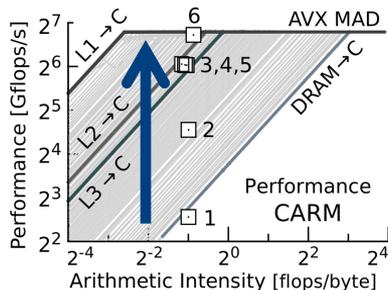|  | **Architectural Efficiency** | $nw = 1$ | $nw = 2$ | $nw = 3$ | $nw = 4$ | $nw = 5$ | $nw = 6$ |
|---|---|---|---|---|---|---|---|
| **FMA** | KNL | 84.98% | 77.50% | 66.77% | 55.28% | 46.56% | 39.65% |
|  | V100 | 97.36% | 91.50% | 76.70% | 65.44% | 65.07% | 66.38% |
|  | **Performance Portability** | **90.76%** | **83.92%** | **71.39%** | **59.93%** | **54.28%** | **49.65%** |
| **No-FMA** | KNL | 82.06% | 72.95% | 73.74% | 78.72% | 81.28% | 82.81% |
|  | V100 | 92.88% | 92.88% | 97.43% | 98.91% | 1 | 99.73% |
|  | **Performance Portability** | **87.14%** | **81.72%** | **83.95%** | **87.67%** | **89.93%** | **90.49%** |

# 4. Energy Roofline

**Performance, Power Consumption, Energy Efficiency**

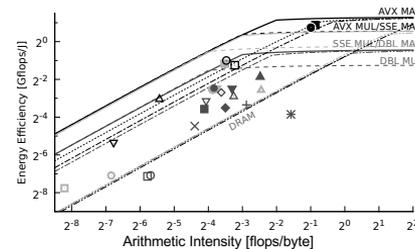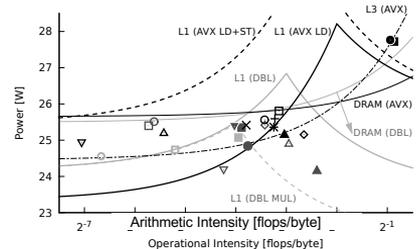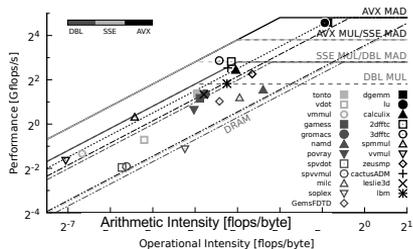# Energy Roofline - GEMM



**Cache-aware Roofline Models**

- **Power Consumption based on CARM**
  - **Relates Watts with FLOPs/bytes**
  - **Defines power envelope for different types of FP and memory operations**

| VERSION | OPTIMIZATION STRATEGY |
|---------|----------------------|
| 1 | Basic implementation: Row-major matrices |
| 2 | Improved memory access by transposing B matrix |
| 3, 4, 5 | Blocking for caches: L3 (pt. 3), L2 (pt. 4) and L1 (pt. 5) |
| 6 | Highly optimized Intel MKL implementation |

- A. Ilic, F. Pratas, and L. Sousa, "Cache-aware Roofline model: Upgrading the loft", IEEE Computer Architecture Letters (2013)
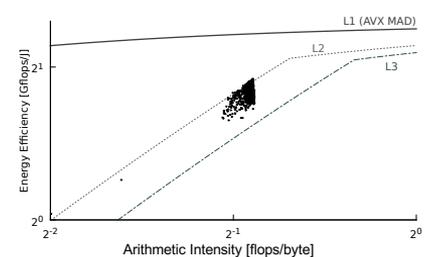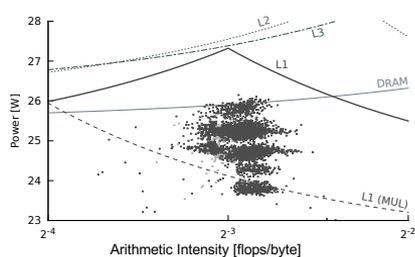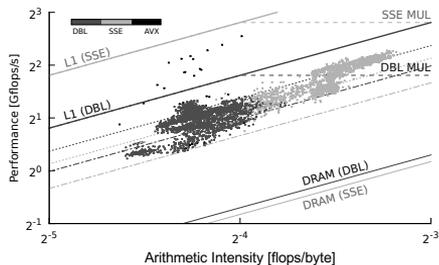- A. Ilic, F. Pratas and L. Sousa, "Beyond the Roofline: Cache-aware Power & Energy-Efficiency...", IEEE Transactions on Computers (2017)

# Energy Roofline - Use Cases



**Application Characterization**

**Online Monitoring**
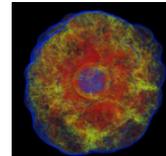
- Ilic, A., Pratas, F. and Sousa, L., "Beyond the Roofline: Cache-aware Power & Energy-Efficiency Modeling…", IEEE Transactions on Computers (2017)
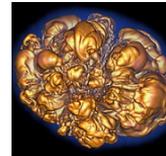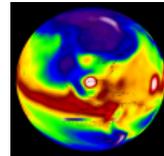- Antão, D., et.al.,"Monitoring Performance and Power for Application Characterization with CARM", PPAM'13
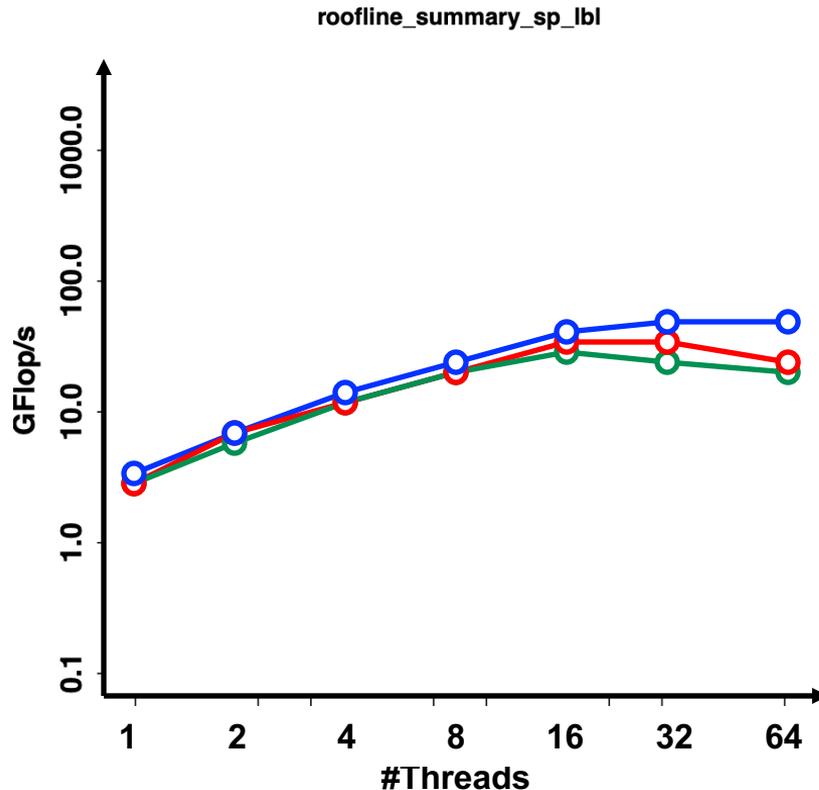
# 5. Scaling Trajectories

**What's causing bad scaling from Roofline point of view?**
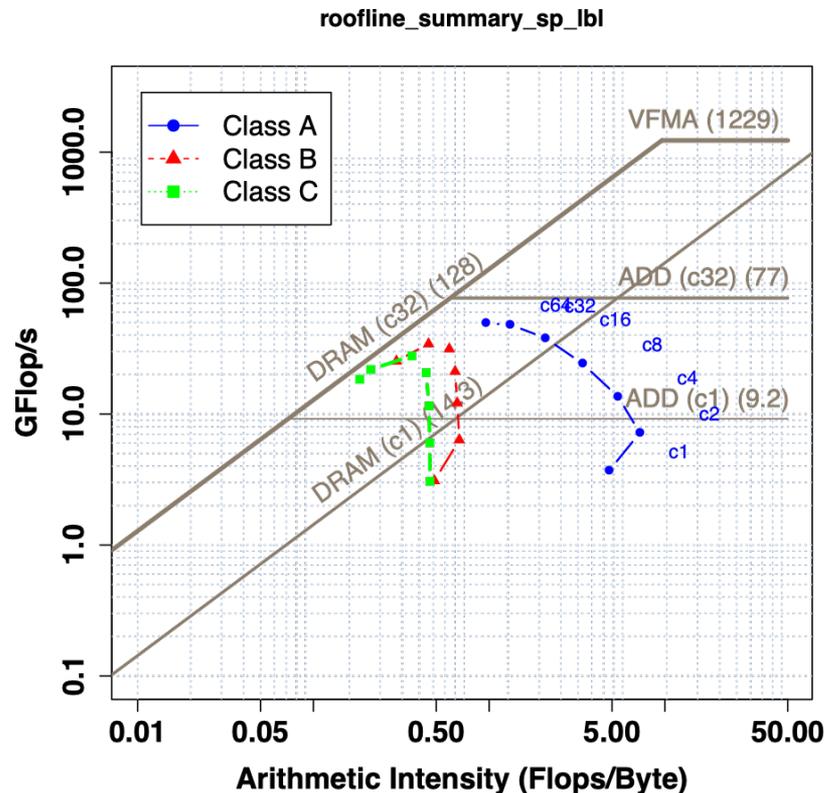
# Roofline Scaling Trajectories

- **We often plot performance as a function of thread concurrency**
  - **Carries no insight or analysis**
  - **Provides no actionable info**

**roofline_summary_sp_lbl**

# Roofline Scaling Trajectories

- ■ We often plot performance as a function of thread concurrency
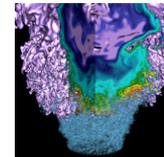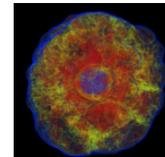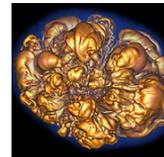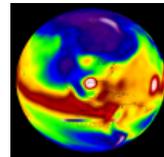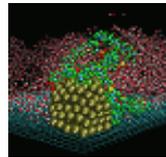  - ○ Carries no insight or analysis
  - ○ Provides no actionable information.

- ■ Use Roofline to analyze thread (or process) scalability
  - ○ 2D scatter plot of performance as a function of intensity and concurrency
  - ○ Identify loss in performance due to increased cache pressure (data movement)



roofline_summary_sp_lbl

• Khaled Ibrahim, Samuel Williams, Leonid Oliker, "Roofline Scaling Trajectories: A Method for Parallel Application and Architectural Performance Analysis", HPBench, July 2018
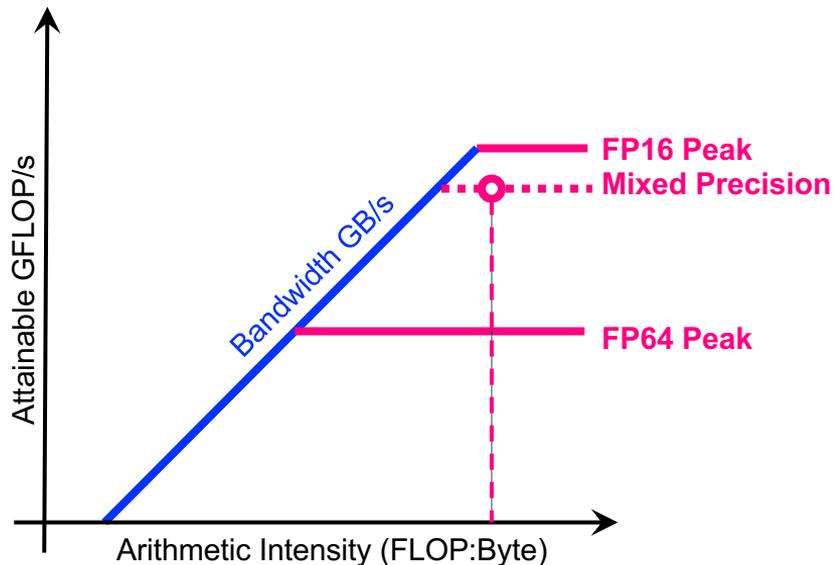
# 6. Mixed Precision

## FP64, FP32, FP16, CPU, GPU

# Mixed Precision

Benefits of reduced/mixed precision:

- From FP64 to FP32
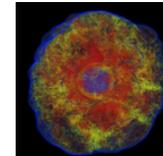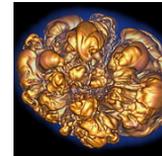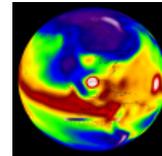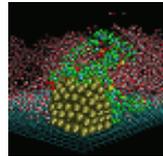  - 2x due to bandwidth savings or compute unit availability
  - similar for network communication
- More support on modern architectures
  - ~15x FP16 over FP64 for some ops

NESAP collaboration with **CRD (Costin Iancu)** and **NVIDIA (Chris Newburn)**
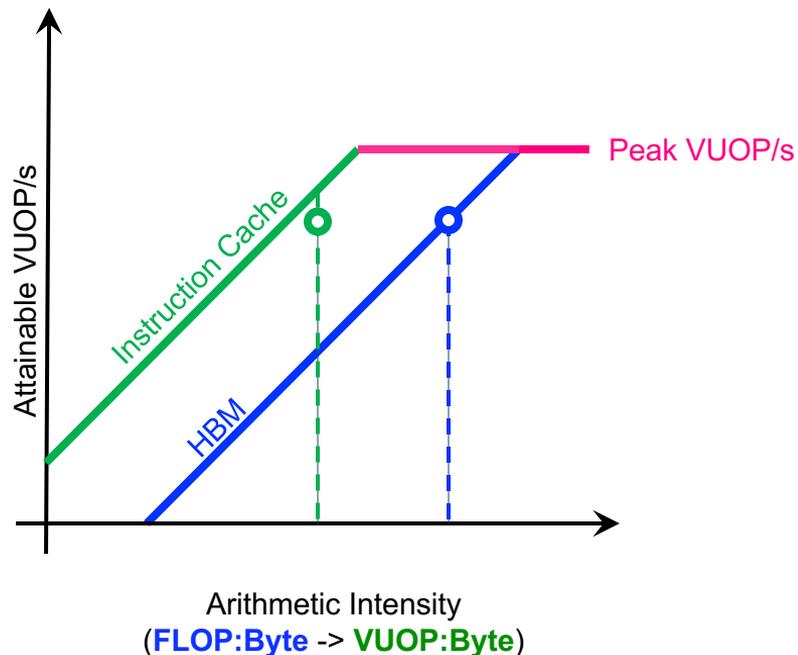
# 7. Instruction Roofline

**FLOP, INTOP, IPC**

# Instruction Roofline

- FP instructions can be the minority in many HPC codes
- Emerging domains have ~no FP
  - Graphs
  - Hash tables
  - Bloom filters
  - Searches
- FLOPs is agnostic of precision, scalar/vectors/tensors, …

- **Instruction Roofline**



Arithmetic Intensity
(**FLOP:Byte** -> **VUOP:Byte**)

# Instruction Roofline

## FLOPs-based Roofline

- FMA doesn't change Arithmetic Intensity (FMA == FMUL+FADD)
- Vectors/tensors don't change Arithmetic Intensity
- Vector integer operations don't change Arithmetic Intensity
- Reducing precision (64b, 32b, 16b) **increases** Arithmetic Intensity
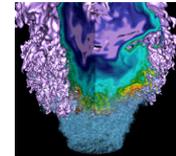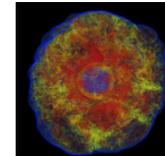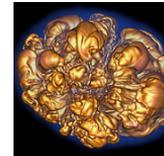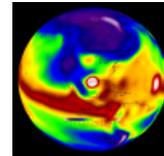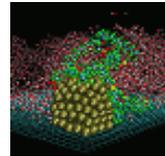
➢ **Tells us about performance**

## VUOPs-based Roofline

- FMA cuts Arithmetic Intensity in **half** (half the number of VUOPS)
- vectors/tensors reduce Arithmetic Intensity (**SIMD cuts VUOPS by 8x**)
- Vector integer operations **increases** Arithmetic Intensity
- Changing precision doesn't change Arithmetic Intensity

➢ **Tells us about VPU/pipeline utilization and bottlenecks**

# 8. Empirical Roofline Toolkit (ERT)

## Machine Characterization, Peak FLOP/s, Bandwidths
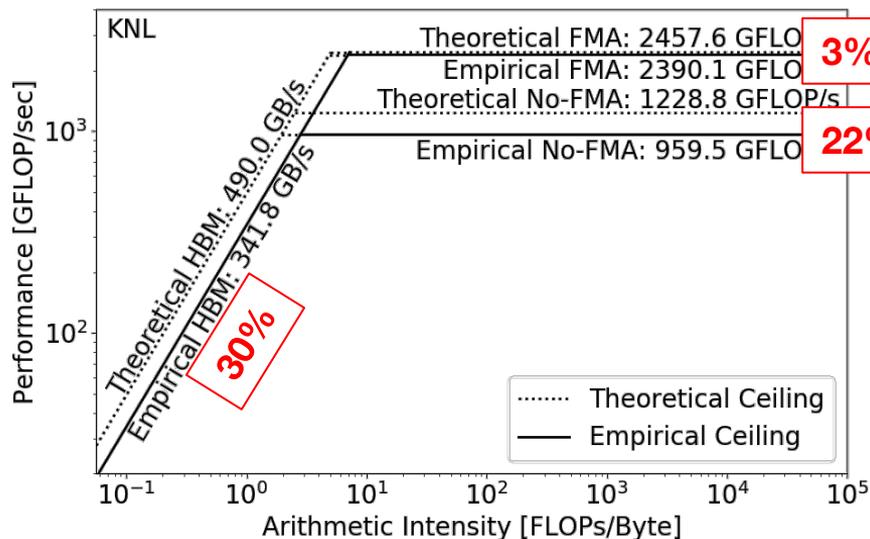
# Empirical Roofline Toolkit (ERT)
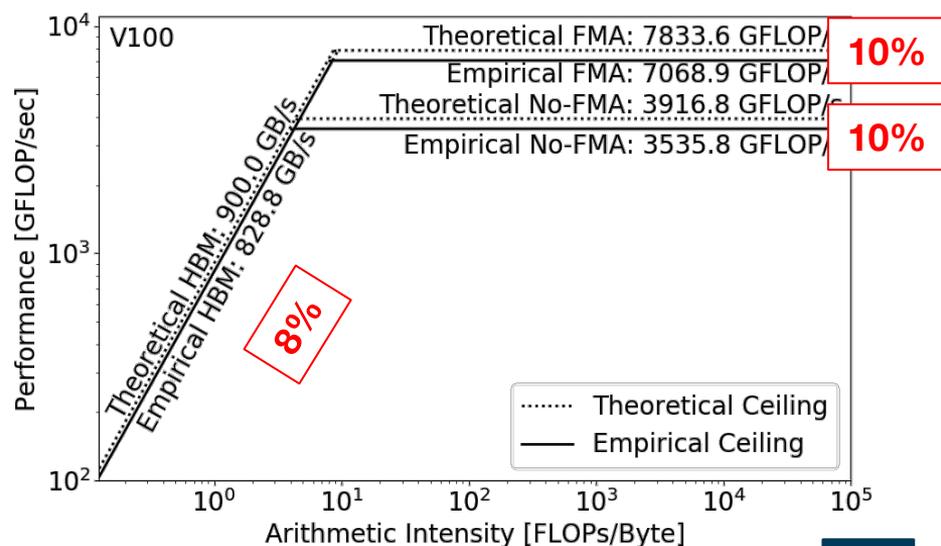
Theoretical compute ceiling on KNL:

$$64 \text{ cores} \times 8 \text{ DP/vector} \times 2 \text{ FLOPs/FMA} \times 2 \text{ vectors} \times 1.2 \text{ GHz} = 2.46 \text{ TFLOP/s}$$

Theoretical compute ceiling on V100:

$$80 \text{ SMs} \times 32 \text{ FP64 cores/SM} \times 2 \text{ FLOPs/FMA} \times 1.53 \text{GHz} = 7.83 \text{ TFLOP/s}$$
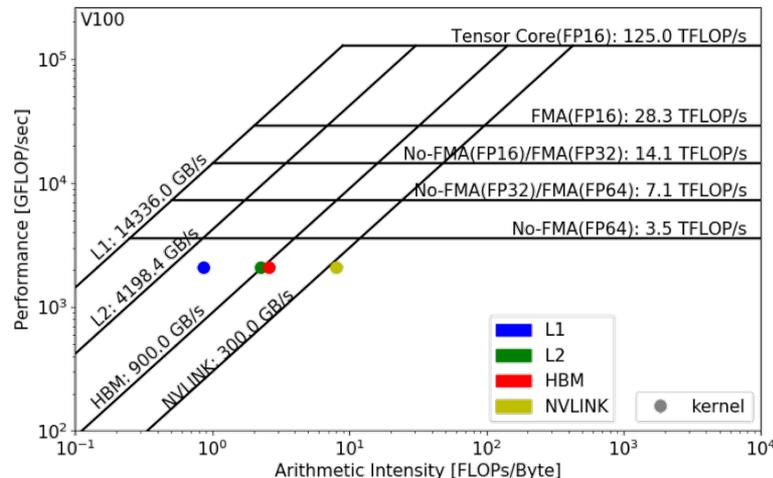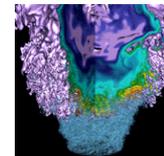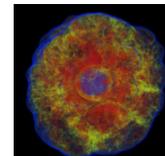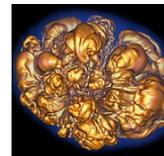


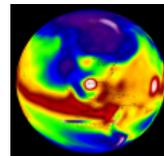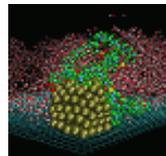Cori KNL partition

Voltar at UOregon

# Empirical Roofline Toolkit (ERT)

- **ERT can't detect all the ceilings yet - IN DEVELOPMENT!**
  - **Haswell/KNL:   L1, L2, L3/HBM, DDR**
  - **V100:        L2, HBM, DDR**

- **Our goal is to incorporate**
  - **the full memory hierarchy**
  - **instruction mix (e.g. FMA/no-FMA)**
  - **data type (e.g. FP64, FP32, FP16)**
  - **compute units**
    **(e.g. CPU/CUDA core/Tensor core)**

- **Ceilings can be omitted if irrelevant**



Empirical Roofline Toolkit (ERT). https://bitbucket.org/berkeleylab/cs-roofline-toolkit/

# Closing

# The Roofline Tree is Flourishing

LBNL CRD Roofline Research:

https://crd.lbl.gov/departments/computer-science/PAR/research/roofline/publications/

**Collaborate with us!**

# Acknowledgement

**Thank You**